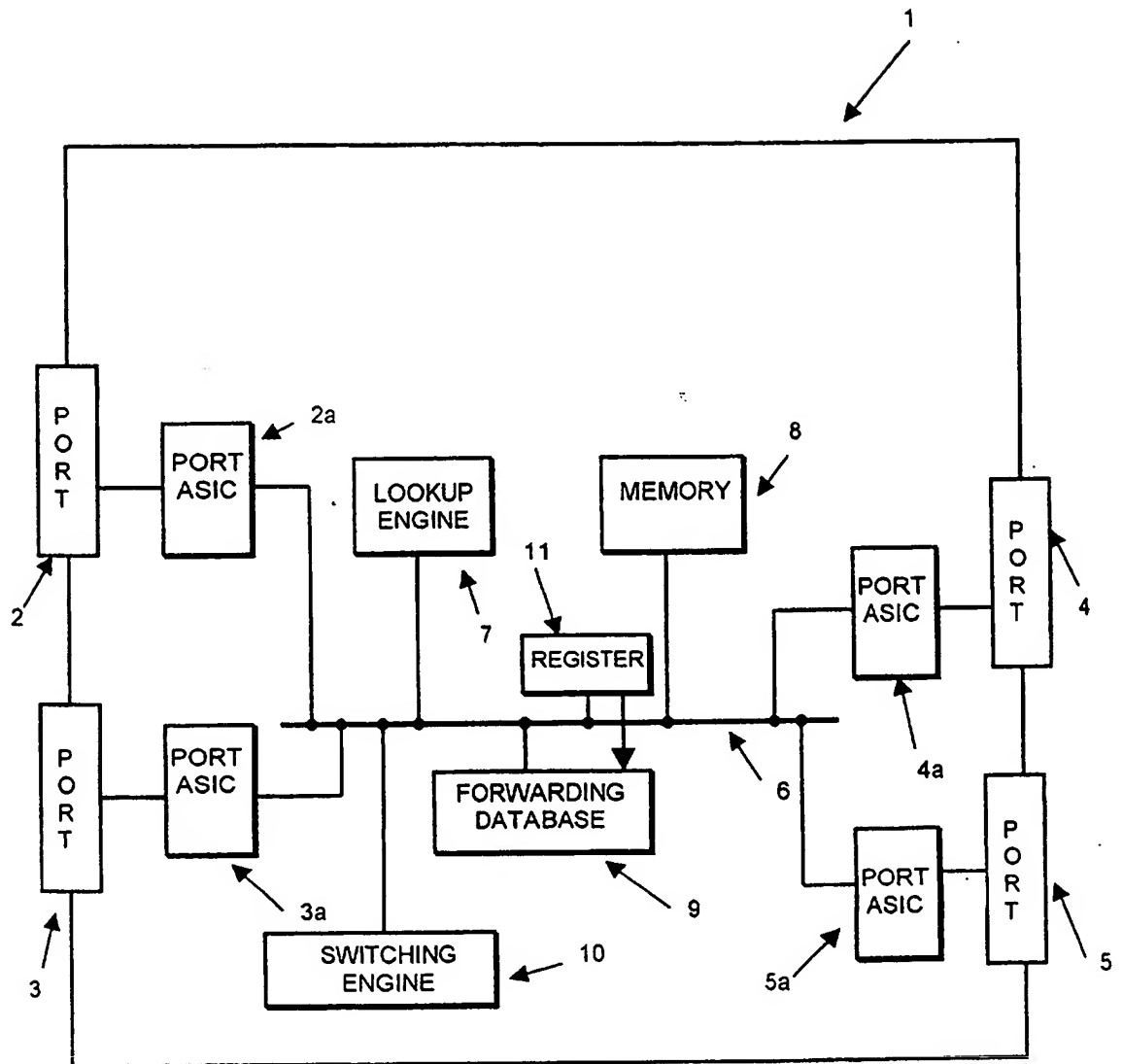


(43) Date of A Publication 01.08.2001

GB 2 358 760 A

FIG.1



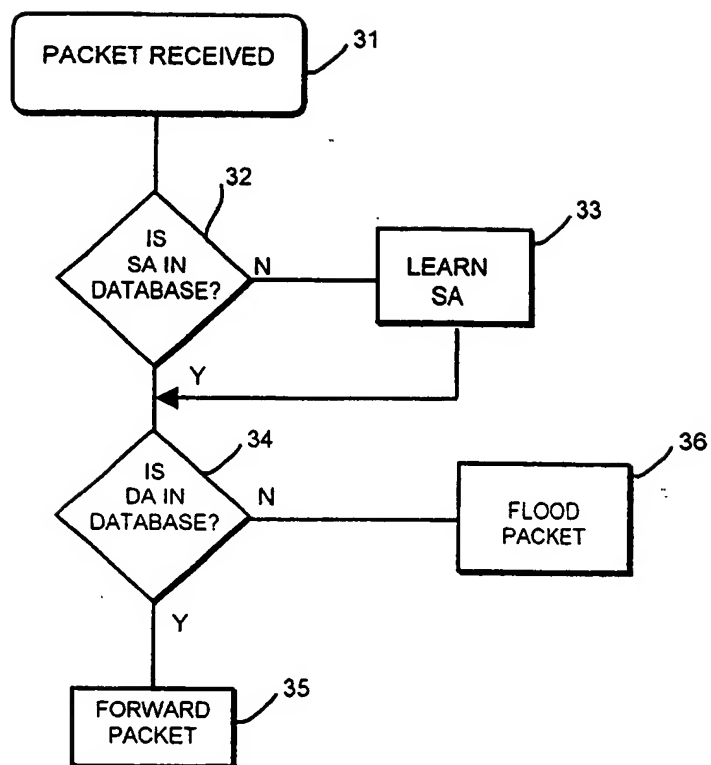


FIG.3

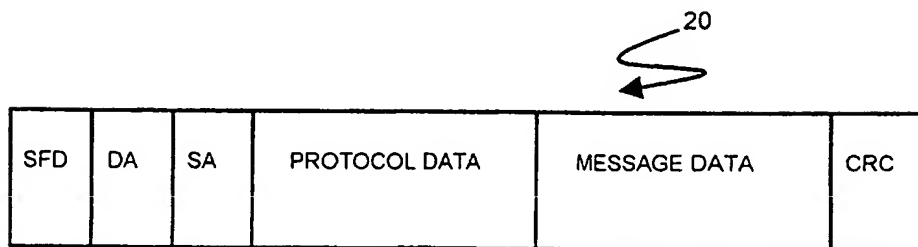


FIG.2

3/7

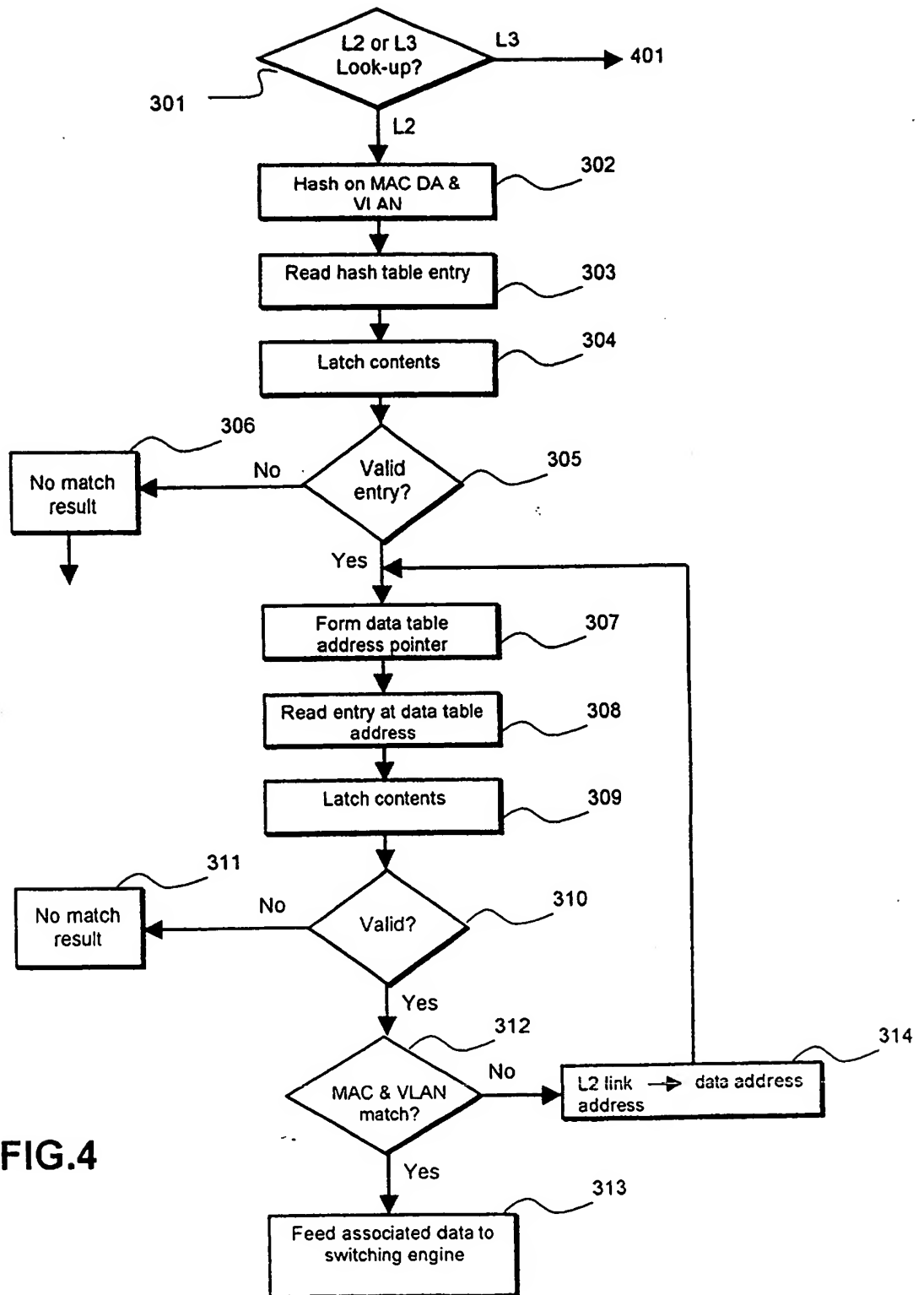


FIG. 4

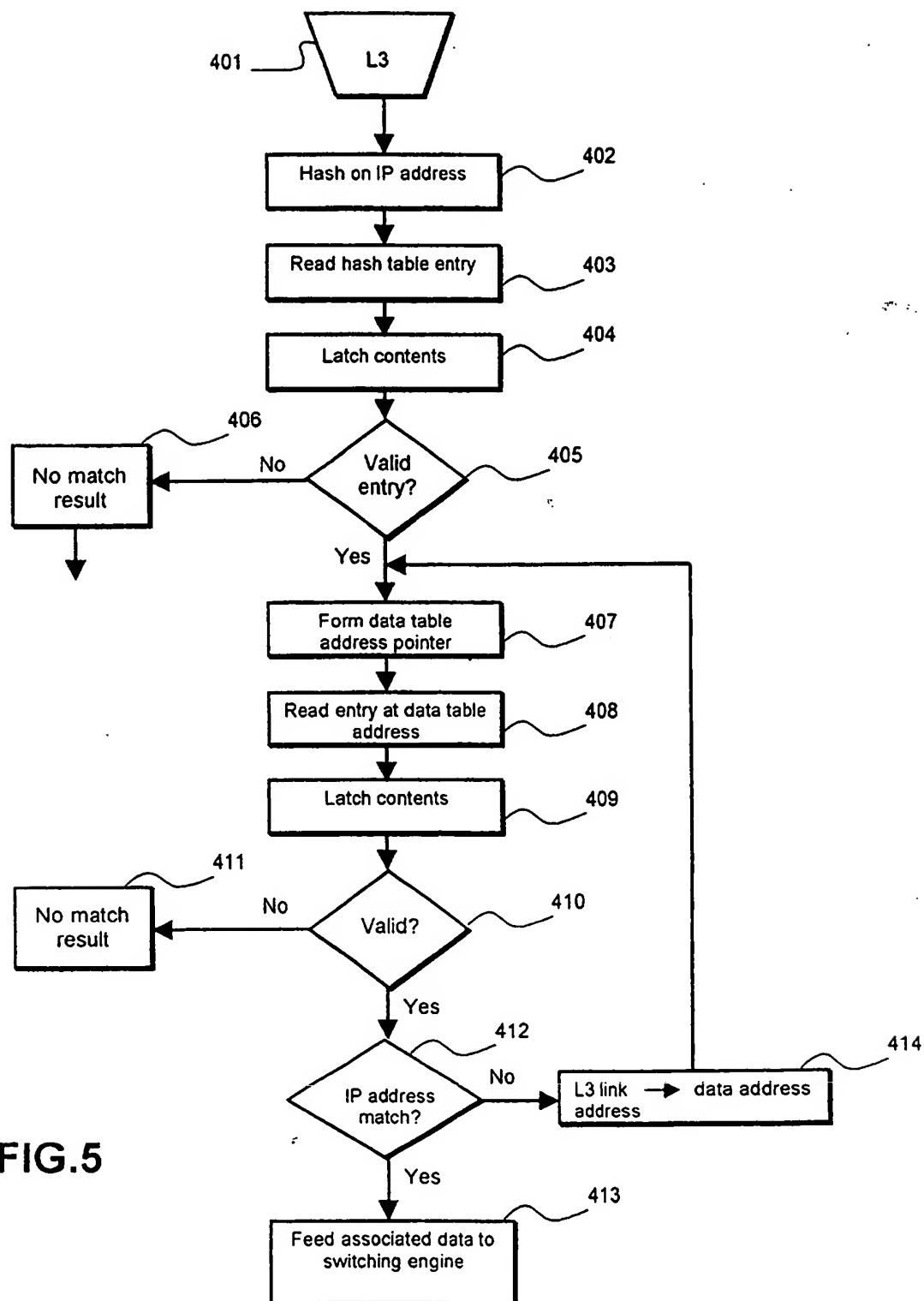


FIG. 5

MAC	IP	VLAN	PortMask	Age

FIG.6

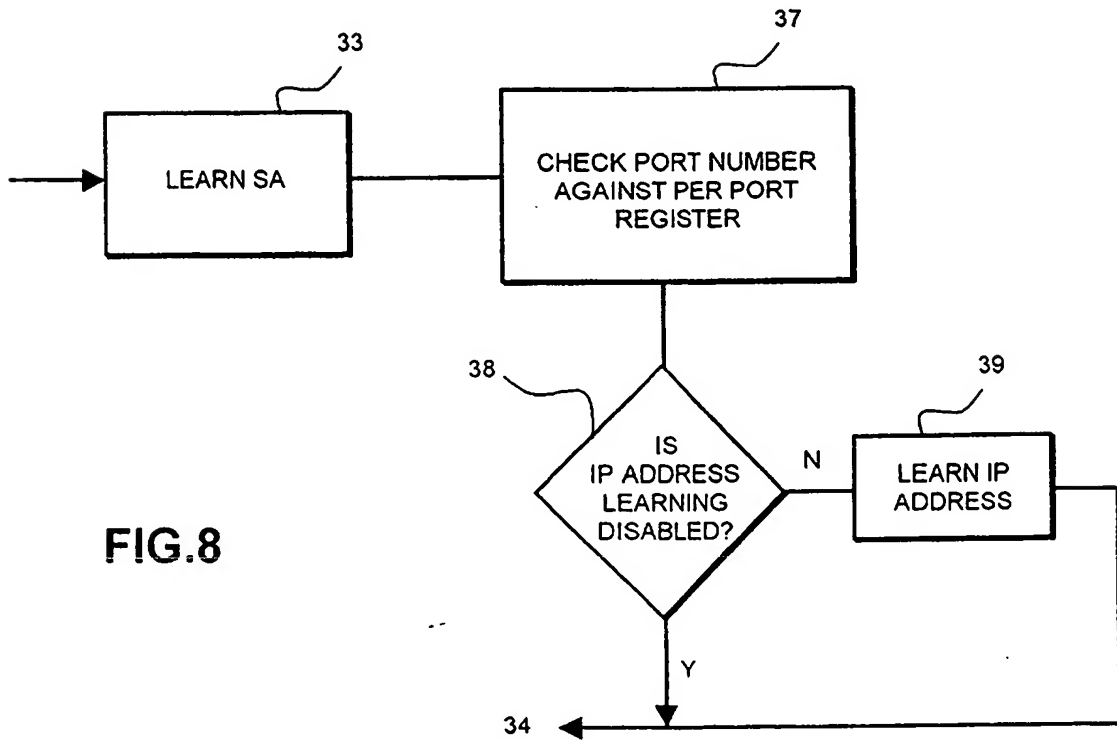


FIG.8

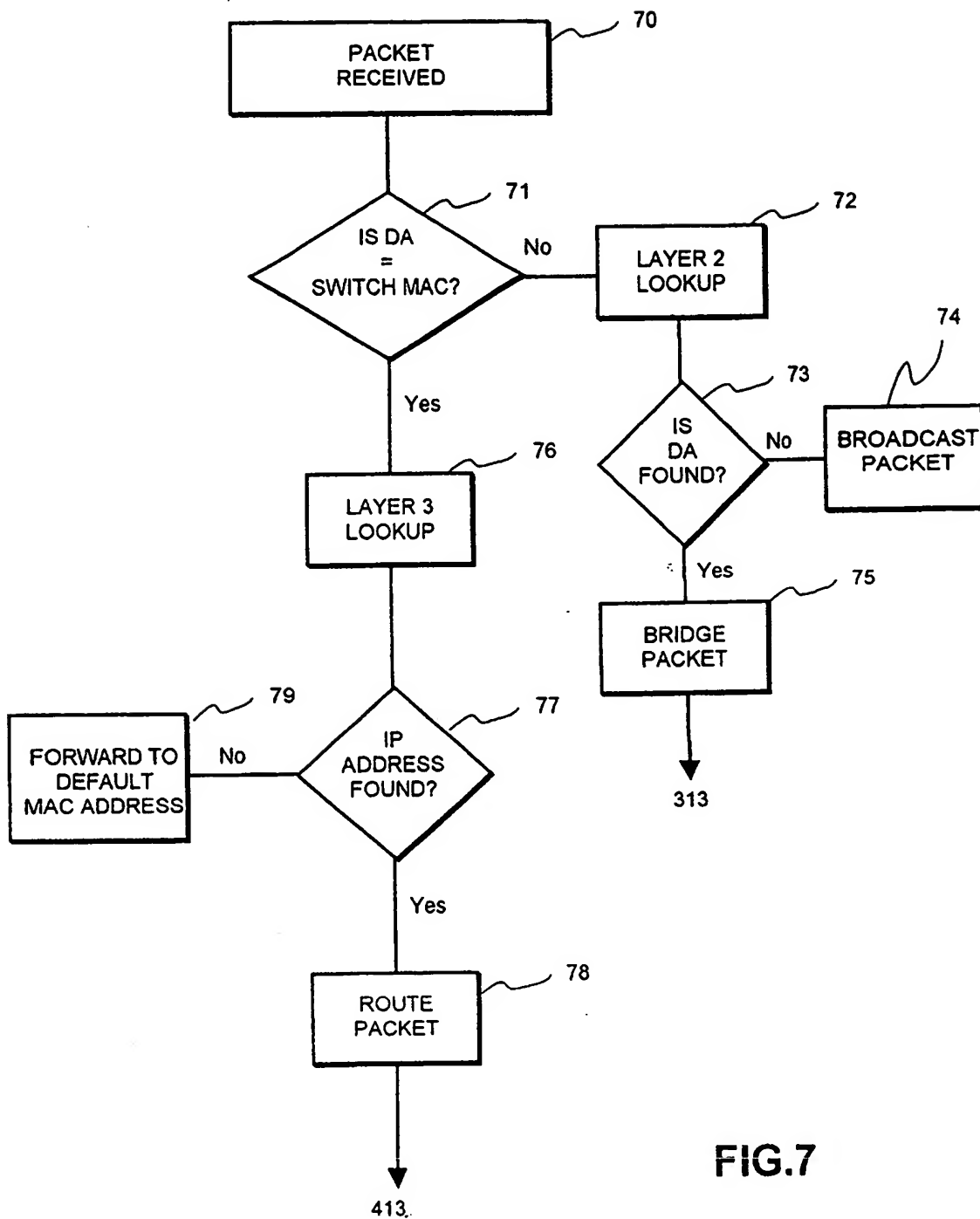
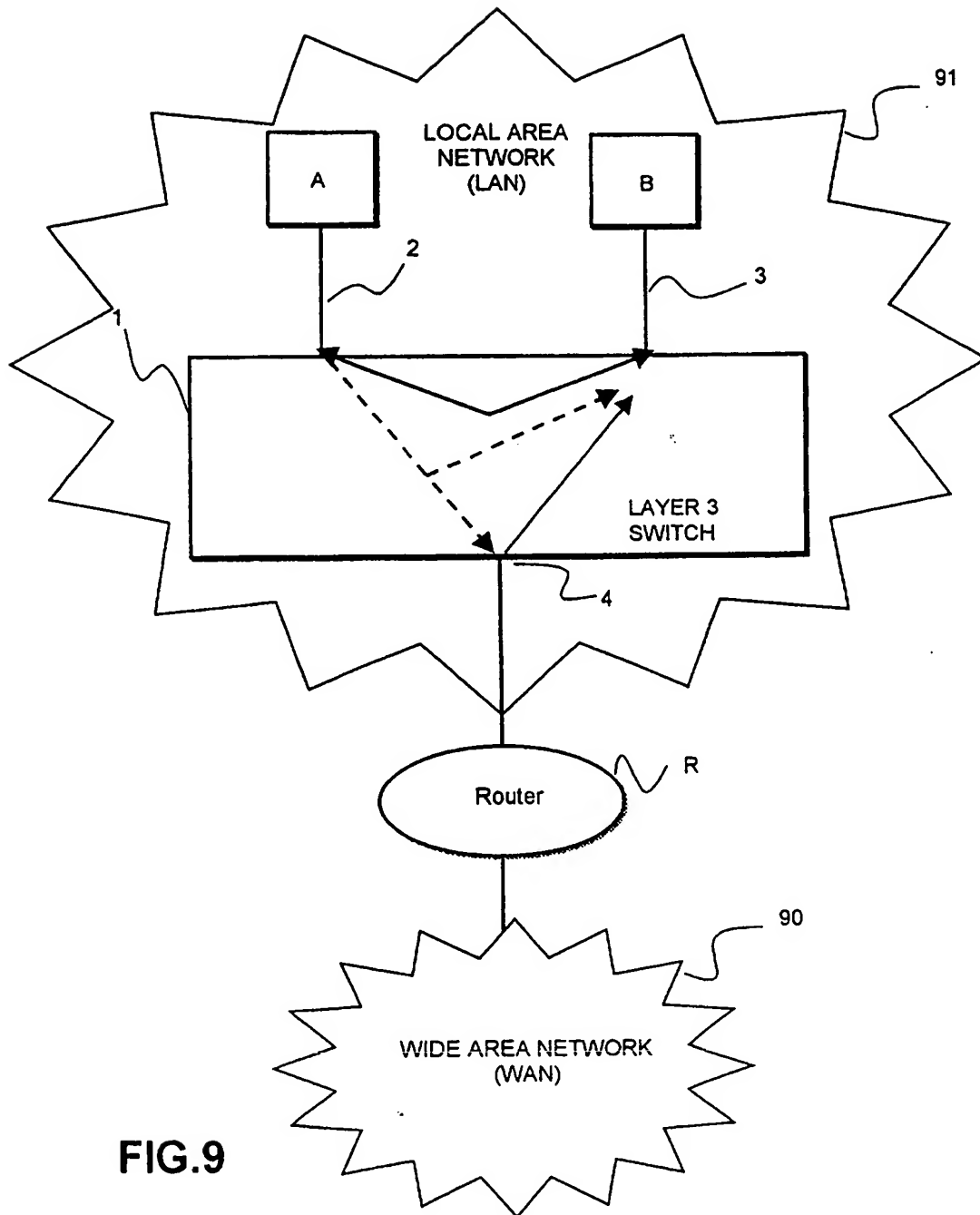


FIG.7

**FIG.9**

NETWORK SWITCH WITH SELF-LEARNING ROUTING FACILITY

Field of the Invention

5 The present invention relates generally to packet-based communication systems wherein data packets including address data and message or control data or both are propagated about a network in accordance with address data in the packets. The invention particularly relates to a network switch which includes a forwarding database and a multiplicity of ports of which one is connected to a router. The invention is mainly intended to facilitate
10 the insertion into a network of a switch which can respond to protocol addresses and be able to relieve the router of traffic which requires simple routing decisions.

Background to the Invention

15 As is well known, a data packet is typically formed in a relatively high level of a communications protocol and before it is transmitted from an originating device it has attached to it a header which includes address data. The address data normally includes a protocol or network address, defining a 'subnet' on which the destination station is located and usually also an identification of that destination station. The network layer or protocol
20 address is employed by a 'router', which term is intended to include devices which have a routing facility, to direct the packet to the appropriate subnet.

The address data within a packet needs to include at some stage a media access control address, otherwise known as a 'layer 2' or 'data link' address. The media access control
25 (MAC) address is employed by switches and other devices to determine, when forwarding a packet, the specific device to which the packet should be sent.

When a packet does not contain a media access control address, as when for example a first end station initially attempts to communicate with a destination end station, it is
30 necessary to perform an address resolution protocol, wherein a packet is broadcast indiscriminately. An end station receiving an address resolution packet (ARP packet)

containing its protocol address can reply with its media access control address. That enables a switch to establish in a forwarding database an entry which relates the particular protocol address with a media access control address and typically a port number of a port to which packets addressed to that destination end station will be sent.

5

Some operations in response to address data in packets are comparatively simple and speedy. For example, 'bridging' is the conventional term employed for responding to the MAC address and directing a packet to the device identified in that address. However, other forms of response, particularly 'routing' are more complex and require greater time. In particular, it is necessary to employ a router to perform such tasks as determining a best route for a packet to take, the prevention of indefinite looping of a packet, and a variety of other functions well known to those skilled in the art. Commonly, the performance of those ancillary functions is very much slower than the performance of a bridging function.

10

However, it is known to provide a switch which can operate both in 'layer 2' and 'layer 3', and which more particularly can in a default mode perform layer 2 look-ups but which can act also as a router, performing layer 3 look-ups. One example of such a switch, which operates with a single look-up table for both bridging and routing decisions, is disclosed in published GB patent application No. 2337674. Another example is a 12 port 100/1000 Mb/s Ethernet switch type 3C17700 made by 3Com Corporation. Such switches require configuration of their routing tables in order to operate in 'layer 3' but differ from fully functional routers in that the routing can be performed at high speed, e.g. 'wire speed'. It is presumed in the following that the 'network switch' employed is a switch of this character, being capable of 'layer 2' switching and, with appropriate configuration, 'layer 3' switching. Such a switch does not perform all the functions which a software controlled router can normally perform. If such a network switch is employed for example in a local area network and a router is also provided to route packets between that local area network and (for example) a wide-area network it would be beneficial to offload from the router the basic routing function (layer 3 switching) within the local area network and in particular, irrespective of what kind of router is employed, to employ the layer 3 switching capability of the switch to route local traffic between subnets.

20

25

30

As will be apparent to those skilled in the art and as more particularly discussed hereinafter, a switch which is capable of 'layer 3' routing will include in its forwarding database entries which relate a protocol (IP) address, a MAC (media access control) address identifying the next hop of a package intended for the protocol address and an
5 identification (such as a number or port mask) of the port to which a packet routed according to that entry must be sent by the switch. The action of establishing in a database an entry of that character (which may include an identification of a virtual local area network) is usually termed 'learning' the protocol address for a particular port.

10 It is known in itself to control a switch, for example remotely by way of a 'management' port, so that protocol (layer 3) addresses cannot be 'learned' for a particular port.

Summary of the Invention

15 The main object of the invention is to facilitate the offloading of routing decisions from a router to an associated switch which is capable of switching decisions on both media access control address and protocol addresses. The main feature of the invention is to provide the connection between the router and the switch only by way of a port in respect of which the switch can learn media access control addresses but is unable (for example
20 by being specifically disabled) to learn protocol (IP) addresses.

Further advantages and features of the invention will become apparent from the following detailed description with reference to the accompanying drawings.

25 Brief Description of the Drawings

Figure 1 is a general schematic representation of a known form of switch which may be employed in the invention.

30 Figure 2 is a simplified illustration of an addressed data packet.

Figure 3 is a flow diagram illustrating a learning process for a network switch.

Figure 4 is a flow diagram of principally a 'layer 2' look-up process in a network switch.

5 Figure 5 is a flow diagram of a 'layer 3' look-up process in a network switch.

Figure 6 is a simplified schematic illustration of a data table.

10 Figure 7 is a flow diagram of the operation of a switch capable of layer 2 and layer 3 switching decisions.

Figure 8 is a partial flow diagram of an address learning process.

15 Figure 9 is a schematic illustration of a connection of a switch and a router in accordance with the invention.

Detailed Description

20 Although the specific instruction of a switch is not necessarily an important feature of the invention, provided that the switch has both the storage ability and the processing ability that the invention requires, Figure 1 is intended to show schematically the basic components of a switch that is suitable for use in the present invention. Typically, switches have twelve or twenty-four ports or even more. For the sake of simplicity, the switch 1 shown in Figure 1 has only four ports, identified as ports 2, 3, 4 and 5. As will be
25 seen later, it will be assumed that ports 2 and 3 are connected to other network devices, port 4 is connected to a 'router' and port 5 is a management port by means of which the switch can be configured by remote control in a manner well known to those skilled in the art.

30 If, as is preferred, the switch 1 is primarily a hardware switch, the various components within the switch 1, apart from most of the memory, be provided on a single ASIC

(application specific integrated circuit). However, for ease of explanation, the various components of the switch are separately shown in Figure 1. In this example therefore, each of the ports 2, 3, 4 and 5 has a respective 'port ASIC', 2a, 3a, 4a and 5a respectively. These components include the media access control devices (MACs) which perform (known) operations on packets entering and leaving the switch while the packets are in a format independent of the particular transmission medium to which a respective port is connected. The port ASICs also include a 'physical layer device' which not only converts packets from a media independent format to a format appropriate for the particular transmission medium but also includes various other functions such as for example auto-negotiation, particularly in the case of 'Ethernet' networks described in IEEE Standard 802.3.

The switch 1 includes a bus system 3 by means of which packet data and control and status data are conveyed between the various components of the switch. The switch includes a look-up engine 7, the operation of which will be described later, a memory 8 which may be employed for the temporary storage of packets in 'queues' before they are sent to their destination ports, a forwarding database 9, which will be described with reference to Figure 6, and a switching engine 10. The switching engine will retrieve packets temporarily stored in memory 8 and direct them to respective ports in accordance with, for example, a port mask obtained from a relevant entry in the forwarding database 9. The switch also includes a register 11 the function of which will be explained later.

Figure 2 illustrates in simplified schematic form a typical packet employed for the conveyance of data in a packet-based data communication system in which a switch such as switch 1 may form part. The packet comprises a start-of-frame delimiter (SFD), media access control address information, comprising a destination address (DA) and a source address (SA), protocol data, message data and cyclic redundancy check data (CRC). The media access control addresses define, if they are present, the source and destination devices in one 'hop' of a packet. The protocol data includes network address data defining, for example, the network to which the ultimate destination of the packet belongs

and usually also an identification of a device within that network. The message data need not be present, as in the case of a control packet.

5 Figure 3 illustrates mostly the learning process for MAC addresses typical of a network switch. A packet is received, stage 31, and a look-up, performed by means of look-up engine 7 in forwarding database 9, determines whether the source address (SA) is already the subject of an entry in the database. If it is not, then the address is 'learned' (stage 33), that is to say made the subject of a table entry including an identification of the port on which the packet was received and a VLAN number. If the switch is to be used for routing 10 (layer 3 switching) as well as bridging (layer 2 switching), an entry will typically include the protocol (IP) address of the packet.

15 In ordinary, layer2/layer switches, IP addresses may be learned at this stage. The switch includes on its ASIC a per port register 11 which identifies those ports for which IP addresses may not be learned. This will be further explained with reference to Figure 8.

20 In order to determine where the packet should be sent, a further look-up is made (stage 34) to find a match for the destination address (DA) in the database. If the address is found, then the packet may be forwarded (stage 35) from the port associated with that MAC address in the forwarding database. For this purpose the entry is read out from the forwarding database and fed to the switching engine 10.

25 If it should happen that the destination MAC address is not in the forwarding database, it is normally necessary to 'flood' or 'broadcast' the packet (stage 36). By this is meant that a copy of the packet is supplied to all (or all of a selected plurality) of the ports in order to obtain an ARP (address resolution protocol) response from a device having the network address identified in the packet. That device will respond with its MAC address and enable this address to be learned in respect of the relevant port in the forwarding database.

30 Figure 3 (particularly stage 34) is intended to include the case when the MAC destination address (DA) of the packet matches the MAC address of the switch; if the packet is of

appropriate IP type, it can be routed (stage 35). If the destination IP address is not in the database the packet would be sent by a default route, but not broadcast as in the case of an (unsuccessful) layer 2 look-up.

5 Figures 4 and 5 will be discussed in relation to Figure 6, which illustrates a typical 'combined' data table which by way of example may perform the functions of a 'routing table', a 'bridging table' and an 'ARP cache'. This shows a forwarding database which contain a multiplicity of entries which may contain a MAC address, a subnet or VLAN address, a network (IP) address, a port mask and an age field. This database is accordingly
10 organised as described in the aforementioned GB-A-2337674.

Figures 4 and 5 illustrate the manner of performing look-ups in a forwarding database. These Figures correspond to Figures in the aforementioned GB-A-2337674. Both Figures assume that the process of look-up is facilitated by means of the hashing of the address
15 which is the subject of the look-up (whether this be a combination of the destination address and VLAN address or the IP address).

Referring first to Figure 4, stage 301 illustrates a decision stage determining whether the MAC address is within a local range of MAC addresses. If the MAC address (DA) is
20 within that local range but the packet type is not IP, then the packet must be bridged. The decision process associated with decision stage 301 will be described with reference to Figure 7.

If a layer 2 look-up is to be performed, the switch will perform a hash operation on a
25 combination of the packets destination address (DA) and VLAN number, stage 302, the hash table entry is read (stage 303) and the contents latch (stage 304). The next stage (305) is an examination whether the entry is valid. This need mean no more than a determination whether the entry is still current or has been aged. If the entry is not valid, then the search fails (stage 306). If the entry is valid, then a data table address pointer is formed from the
30 latched contents of the hash table entry, stage 307, the entry is read, stage 308, the contents of the entry are latched, stage 309, and a determination whether the entry is valid

is made (stage 310). If the entry is not valid, then no result has been obtained, stage 311. If the entry is valid then it is determined, stage 312, whether the MAC address and VLAN number in the entry match those of the destination address which is being looked up. If they do not, then it may be necessary to search another entry, linked to the first by means of a pointer. This expedient is necessary because the hashing of addresses may mean that a plurality of addresses may hash to the same entry. However, both the hashing of addresses, and the use of link pointers are merely preferred features of the switch described in the aforementioned prior application and are not essential to the present invention. The important matter is whether a look-up has found the destination address in the forwarding database. If it has been, the response to stage 312 being 'yes' then the data associated with the entry (such as the port mask) are fed to the switching engine, stage 313.

Figure 5 illustrates a similar look-up which may be performed in respect of layer 3 addresses. In this case, the entry stage is stage 401. The IP address is hashed (402), the hash table entry is read (403) and the contents thereof latched (404). On examination (stage 405) of the validity of the entry, there is no match result if the table entry is invalid (406). If the table entry is valid then a data table address is pointer is formed (stage 407), the entry at the data table address is read (stage 408), the contents latched (409) and a test of validity made (stage 410). If the table entry is not valid then there is no match result (stage 411). If the table entry is valid then there is a test to determine whether the IP addresses match, stage 412. If there is a list of addresses linked by pointers, because of the use of hash tables, then a link pointer 414 points to another address in the table and the loop from stage 414 via 407 to stage 412 is reiterated.

Again, however, the important matter is whether a match of the IP address in the packet has been found with an address in the data table. If so, then the relevant data in that entry including the port mask is fed to the switching engine, shown in stage 413.

Figure 7 is a summary of the decision process in stage 301 as well as a summary of the layer 2 and layer 3 look-ups shown in Figures 4 and 5. By way of introduction to Figure 7,

it should be remarked that the switch 1 is normally, for example by way of the management port 5, configured with its own MAC address.

5 Referring now specifically to Figure 7, if a valid IP packet, that is to say a packet containing a valid 'network' or layer 3 address, stage 70, a determination (stage 71) is made to

10 If the incoming valid IP packet does not contain the local MAC address, then the response to stage 71 (which corresponds to stage 301 in Figure 4) is negative and the switch will perform a layer 2 look-up, summarised by stage 72 and more particularly illustrated in Figure 4. If the destination MAC address is found then the packet may be 'bridged', that is to say switched to the relevant port on the basis of the MAC address and port number. If the MAC address is not found in the database then the packet will be broadcast, stage 74.

15 If however, the incoming packet does have a destination MAC address (DA) corresponding to the address of the switch, a layer 3 look-up will be performed (stage 76 and Figure 5). If the layer 3 (network or IP) address is found by for example the process shown in Figure 5, the search will retrieve the next hop MAC address and therefore the relevant port number and the packet can be 'routed'. It will, in well known manner not requiring description, have its MAC source address changed to the MAC address of the
20 switch 1 and the 'TTL' will be decremented.

Figure 7 shows a further situation, where the IP address search, stage 76, yields a negative result. In this case, the packet needs to be forwarded by a default route. This is defined by
25 a 'default' MAC address.

The foregoing description is intended to provide the reader with a background for understanding the learning and look-up processes in a layer 2/layer 3 switch of the kind which may form a combination with a router according to the invention.

30

Disabling of Learning

It is known, for example, in a 'local office interconnect' scheme to modify the operation of a switch such as switch 1 by preventing the learning of IP addresses in respect of a selected port.

5

As indicated below the layer 3 switch router listens on the network for 'router alive' messages. When it detects such a message it reads the source port number that the source MAC address has been learned against. It then writes that port number to a per port register which disable learning for that port. It does not disable MAC address learning so

10

'Router alive' messages are transmitted at various intervals and can be detected by snooping on, for example OSPF, PIM and RIP packets.

15

It is also feasible to disable learning for a particular port by local (e.g. manual) programming of the switch.

Whether the learning is disabled automatically or not, the ability to learn IP addresses against a port is controlled by the 'per port' register 11.

20

The learning of IP addresses is illustrated in Figure 8, which is intended to be read in conjunction with Figure 3. Thus, although the switch 1 learns MAC addresses routinely (stage 33), the learning of the IP address of a packet depends on a check (stage 37) of the port number of the incoming packet against the per port register 11. If IP address learning is disabled for that port (stage 38) the learning process reverts to stage 34 of Figure 3. This does not preclude layer 3 switching (routing) if the MAC DA of the packet matches the MAC address of the switch. If IP address learning for that port is not disabled, the IP address is learnt against the respective port (stage 39).

25

30

Offloading of a Router

Reference will now be made to Figure 9 which shows a layer 2/layer 3 switch 1 as described in the foregoing connected by way of port 2 to a first network device A, connected by way of port 3 to second network device B and by way of port 4 only to a router R. As described in the foregoing, the switch 1 is configured so that it is unable to learn IP addresses in respect of port 4.

In this typical example, the router R is a 'standard' router which performs, principally under software control, wide area routing functions for a wide-area network 90. The switch 1 is principally intended for operation on a local-area network (LAN) 91 to route traffic between sub-nets of the LAN. Large routers such as router R must perform considerably more packet analysis than a switch 1 and not all their functions can (unlike switch 1) be implemented in hardware (i.e. in an ASIC). The router R may route traffic between the LAN 91 and WAN 90 but should preferably not route traffic within the LAN (e.g. between stations A and B).

In a first phase, it will be assumed that the switch 1 receives an ARP packet from station A. The switch will recognise such a packet as not an ordinary IP packet, since ARP packets have a different type field within the packet. The switch 1 will forward such an ARP packet to every possible destination, including station B by way of port 3 and the router by way of port 4. It will also be assumed that an ARP response is generated by station B. This response will have the MAC address of station B and this response will be sent by way of switch 1 to the router R. The router will forward the ARP response packet back to station A together with the MAC address of the router.

At this point no 'layer 3' IP addresses will have been learned by the switch 1 though layer 2 addresses will have been learned through standard bridging rules. In particular, the router's source MAC address will have been learned in forwarding database 9 when the router replies by sending the ARP response packet back to station A.

In a second phase, an IP packet is sent by station A. The switch will perform a look-up in respect of this IP packet. On the assumption that the IP source address of station A is not

5
10
15
20
25
30

5
10
15
20
25
30

10
15
20
25
30

15
20
25
30

25

30

Claims

5

10

15

20

25

30

30

said router being connected to the switch only by way of the specific port.

3. A method according to claim 2 wherein said switch is disposed to route data packets
5 within a local-area network and said router is disposed to route data packets within a wide-
area network.

10

15

20

25

30



INVESTOR IN PEOPLE

Application No: GB 0001654.3
Claims searched: 1-3

Examiner: Richard Howe
Date of search: 28 June 2000

Patents Act 1977
Search Report under Section 17

Databases searched:

UK Patent Office collections, including GB, EP, WO & US patent specifications, in:

UK Cl (Ed.R): H4K (KTK) ; H4P (PPS)

Int Cl (Ed.7): H04L (12/56) ; H04Q (11/04)

Other: Online : wpi ; epodoc ; japio

Documents considered to be relevant:

Category	Identity of document and relevant passage	Relevant to claims
A	GB 2 337 674 A (3Com) - see whole document	
A	GB 2 337 659 A (3Com) - see whole document	
A	EP 0 835 009 A2 (Kabushiki Kaisha Toshiba) - see whole document	

X	Document indicating lack of novelty or inventive step	A	Document indicating technological background and/or state of the art.
Y	Document indicating lack of inventive step if combined with one or more other documents of same category.	P	Document published on or after the declared priority date but before the filing date of this invention.
		E	Patent document published on or after, but with priority date earlier than, the filing date of this application.
&	Member of the same patent family		